

Tipo de Regresor	Modelo Matemático	Función de Costo	Estrategia de Optimización
Lineal Regresor	Aproxima una relación entre un conjunto de variables de entrada $X \in \mathbb{R}^{N \times p}$ y una variable objetivo $t \in \mathbb{R}^N$: $y(X, \omega) = \omega^T X + \omega_0$ $y = X\omega + b$.	Error cuadrático medio: MSE $J(\omega) = \frac{1}{2N} \sum_{n=1}^N (t_n - \omega^T x_n)^2$	Solución Analítica Cerrada: $\nabla_{\omega} J(\omega) = -\frac{1}{N} X^T (t - X\omega) = 0$ Solución: $\omega = (X^T X)^{-1} X^T t$ Resuelve el sistema: $X^T X \omega = X^T t$
Lasso	Introduce una regularización L_1 , reduciendo la magnitud de los pesos: $y(X) = \omega^T X + b$.	Minimiza el MSE y mantiene los coeficientes ω pequeños: $J(\omega) = \frac{1}{2N} \sum_{n=1}^N (t_i - \omega^T x_i - b)^2 + \lambda \sum_{j=1}^p \omega_j $ con $\lambda > 0$	Solución numérica Iterativa: → Descenso por coordenadas: On paso a la vez. → Subgradiente Descendente: → Least Angle Regression: Construye la solución paso a paso a medida que ω entra o sale.
ElasticNet	El modelo sigue siendo el mismo $y = \omega^T X + b$, pero agrega regularización L_1 (Lasso) y L_2 (Ridge).	$J(\omega) = \frac{1}{2N} \sum_{n=1}^N (t_i - \omega^T x_i - b)^2 + \lambda \left[\alpha \ \omega\ _1 + \frac{(1-\alpha)}{2} \ \omega\ _2^2 \right]$ $\alpha = 1$ (Lasso), $\alpha = 0$ (Ridge) $0 < \alpha < 1$ (ElasticNet).	Optimización iterativa convexa → Descenso por coordenadas: Parámetro ω_j a la vez $\omega_j \leftarrow \frac{SC_{2j} \lambda \alpha}{1 + \lambda(1-\alpha)}$
Kernel Ridge	Extiende la regresión lineal al espacio de características de un kernel, se proyectan por una función no lineal: $y(X) = \omega^T \phi(X)$. $K(X, X^T) = \phi(X)^T \phi(X')$	Error cuadrático con L_2 : $J(\omega) = \frac{1}{2N} \sum_{n=1}^N (t_i - \omega^T \phi(x_i))^2 + \frac{\lambda}{2} \ \omega\ ^2$. Derivando ω : $\omega = \phi^T (\phi \phi^T + \lambda I)^{-1} t$.	El problema es convexo y diferenciable, se puede resolver de forma analítica: $\alpha = (K + \lambda I)^{-1} t$ Se usa SVD para invertir $K + \lambda I$ establemente.
SGD Regresor = Stochastic Gradient Descent	Mismo modelo de: $y(X) = \omega^T X + b$.	Diferentes funciones de pérdida según la regresión; más común: squared loss: $J(\omega) = \frac{1}{2N} \sum_{n=1}^N (t_i - (\omega^T x_i + b))^2$ Puede incluir L_1 y L_2 y ElasticNet	No calcula el gradiente usando todos los datos, sino una sola muestra por iteración. La solución es iterativa.

Tipo de Regresor	Modelo Matemático	Función de costo	Estrategia de Optimización.
Bayesian Ridge	parte del mismo: $y(x) = w^T x + b$. pero no estima un valor para w y b , sino distribuciones para cada uno.	$\log p(w t) = -\frac{\beta}{2} \ t - xw\ ^2 - \frac{\alpha}{2} \ w\ ^2$ α y β no son fijos, son variables aleatorias.	La inferencia se hace de forma analítica porque las distribuciones son gaussianas.
Gaussian Process Regressor	Es una extensión no paramétrica de la regresión bayesiana lineal. $f(x) \sim gP(m(x), k(x, x))$.	Se ajustan los hiperparámetros del kernel y del ruido: $\log p(t x) = -1/2 t^T (K + \sigma^2 I)^{-1} t - 1/2 \log [K + \sigma^2 I] - N/2 \log (2\pi)$.	2 Etapas: a) Inferencia: Distribución posterior: $p(f_* x, t, x) = \mathcal{N}(f, \text{var}(f))$. $f_* = K^T x (K + \sigma^2 I)^{-1} t$. b) Optimización de hiperparámetros mediante gradiente descendente.
Support Vector Machine Regressor	Busca una función plana $f(x)$ que se aleje de los datos en menos de un margen, penalizando errores que se le alejen: $f(x) = w^T \phi(x) + b$.	El problema de optimización es: $\min \frac{1}{2} \ w\ ^2 + C \sum (e_i + e_i^*)$. $E \rightarrow$ errores que exceden el margen.	Se usan multiplicadores de Lagrange: $\max -1/2 (a - a^*)^T K (a - a^*) - E \sum (a_i + a_i^*) + \sum t_i (a_i - a_i^*)$ Se obtiene la solución usando métodos de optimización cuadrática.
Random forest Regressor	Es un modelo de ensamblado formado por B árboles de decisión independiente: $f(x) = \frac{1}{B} \sum f_b(x)$ $B \rightarrow$ $b=1$ Cada árbol predice un valor y se toma la media.	Cada árbol se entrena minimizando el MSE: $= \frac{1}{N} \sum_{i=1}^N (t_i - y_i)^2$. Se elige la característica y el punto de corte: $\text{split}(j, S) = \arg \min [MSE_{\text{left}} + MSE_{\text{right}}]$.	Tiene 2 niveles de aleatoriedad controlada: 1. Bootstrap Sampling. 2. Feature Bagging. 3. Crecimiento de árboles. 4. Agregación (media final).

3.

Tipo de Regresor	Modelo Matemático	Función de Costo	Estrategia de Optimización
Gradiente Boosting Regressor	<p>El modelo final es una suma ponderada de árboles débiles:</p> $f(x) = \sum_{m=1}^M v_m \cdot h_m(x)$ <p>$M \rightarrow$ Número de iteraciones. $h_m(x)$: Árbol de decisión entrenado. $v \in [0, 1]$: learning rate. $f(x)$: predicción final</p>	<p>Trabaja con cualquier función de pérdida diferenciable $L(y, f(x))$. En regresión clásica: $L(y, f(x)) = \frac{1}{2} (y - f(x))^2$. Minimiza la pérdida total sobre todos los datos.</p>	<p>Aplica el descenso del gradiente, pero en el espacio de las funciones, no en el de los parámetros.</p> <ol style="list-style-type: none"> 1. Función que minimice. 2. Gradientes negativos. 3. Ajusto del árbol al residuo. 4. Paso óptimo del árbol. 5. Actualiza modelo. 6. Repetir M iteraciones.
XG Boost	<p>Es un modelo aditivo como el GBR, pero añade regularización explícita y una aproximación de 2do orden al gradiente.</p> $f(x) = \sum_{m=1}^M f_m(x), f_m \in F$ <p>$F \rightarrow$ conjunto de árboles. $M \rightarrow$ Número de iteraciones.</p>	<p>En cada iteración, minimiza</p> $Q = \sum_{i=1}^n L(y_i, \hat{y}_i) + \sum_{m=1}^M \Omega(f_m)$ <p>$L(y_i, \hat{y}_i)$ es la función de pérdida. $\Omega \rightarrow$ es la penalización sobre la complejidad del árbol.</p>	<p>Utiliza un desarrollo de segundo orden de la pérdida:</p> $L(y_i, \hat{y}_i + f_m(x_i)) \approx L(y_i, \hat{y}_i) + g f_m(x_i) + \frac{1}{2} h f_m(x_i)^2$ <p>$g \rightarrow$ Gradiente. $h \rightarrow$ Hessiano.</p>

Tipo de Regresor	Minimos Cuadrados	Maxima Verosimilitud	Maximo a Posteriori	Bayesiano lineal	Regresión Rígida Kernel	Procesos Gaussianos	Escalabilidad
Linear Regressor	Minimizar el MSE entre predicciones y valores reales. Formulación: $\min_j J(w) = \frac{1}{2N} \sum_{n=1}^N (t_n - w^T x_n - b)^2$ Solución: $w = (X^T X)^{-1} X^T t$ Resuelto a través de: SUB o QR.	Estimador bajo ruido gaussiano: $t_n = w^T x_n + \epsilon_n$ Formulación: $p(t x, w, \sigma^2) = \mathcal{N}(t_n w^T x_n, \sigma^2)$ Relación: $\arg \max p(t x, w) = \arg \min \ t - Xw\ ^2$	Introduce incertidumbre. Se convierte en "ridge regressor" agregando regularización. Solución: $w_{MAP} = (X^T X + \lambda I)^{-1} X^T t$	Considera toda la distribución sobre los pesos. Solución: $p(w t) = \mathcal{N}(w M_N, S_N)$ $p(w) = \mathcal{N}(0, \alpha^{-1} I)$	Espacio de características no lineal. Formulación = $\min_w \frac{1}{2} \ t - \phi w\ ^2 + \frac{\lambda}{2} \ w\ ^2$. Para linear regressor: caso kernel lineal ($K = XX^T$).	Modelar función como distribución gaussiana: $f(x) \sim g(\theta, K(x, x'))$ El linear regressor usa un GP lineal y sin incertidumbre.	Si hay muchas características el costo crece y se deben usar métodos aproximados.
Lasso	$J_{\text{lasso}}(w) = \frac{1}{2N} \sum_{n=1}^N (t_n - w^T x_n - b)^2 + \lambda \sum_{j=1}^J w_j $	La función de costo puede verse como una verosimilitud max regularizada: $\log p(w t) = \log p(t w) + \log p(w)$ $p(w)$ laplaciano	Introduce un prior laplaciano: $\log p(w t) = -\log p(t w) - \log p(w)$ $p(w) = \frac{\lambda}{2} e^{-\lambda w }$	Prior Laplaciano: de la forma: $p(w) = \mathcal{N}(0, \alpha^{-1} I)$ $p(w) = \alpha \cdot e^{-\lambda w }$ Ya no es una gaussiana cerrada.	Kernel Lasso: $\min_w \frac{1}{2} \ t - \phi w\ ^2 + \lambda \ w\ _1$	El Prior ya no es gaussiano. Si se define un ruido con ruido laplaciano, se puede definir como Lasso..	El modelo debe resolverse con métodos iterativos, lo que afecta el tiempo de procesamiento y memoria con p grande.
ElasticNet	El modelo matemático parte de mínimos cuadrados añadiendo $L_1 + L_2$	El segundo término actúa como una penalización bayesiana híbrida (Laplaciana + Gaussiana).	$-\log p(t w) =$ pérdida del ruido gaussiano. $\log p(w) =$ prior mixto ($L_1 + L_2$).	$p(w) = \frac{e^{-\lambda \alpha \ w\ _1}}{2} + \frac{(1-\alpha)}{2} \ w\ _2^2$	Kernel ElasticNet: $J_{\text{ken}}(w) = \frac{1}{2} \ t - \phi w\ ^2 + \frac{(1-\alpha)}{2} \lambda [\alpha \ w\ _1 + \frac{1}{2} \ w\ ^2]$	Se puede considerar como una mezcla de gaussiano y regularizada.	Tiene mejor estabilidad y convergencia por L_2 pero se puede afectar por α , λ y los datos deben estar normalizados.
Kernel Ridge	Versión no lineal y regularizada. Introduce transformaciones no lineales de las variables. $J(w) = \frac{1}{2N} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 + \lambda / 2 \ w\ ^2$	Espacio de características $\phi(x)$ $t_i = w^T \phi(x_i) + \epsilon$ Agregar regularización L_2 .	$p(t w) = \mathcal{N}(\phi w, \beta^{-1} I) \Leftrightarrow p(w) = \mathcal{N}(0, \alpha^{-1} I)$ El log posterior es: $J(w) = \beta/2 \ t - \phi w\ ^2 + \alpha/2 \ w\ ^2$ $\lambda = \alpha / \beta \rightarrow$ Función de costo	$p(w t) = \mathcal{N}(w M_N, S_N)$ $S_N^{-1} = \alpha I + \beta X^T X$ $M_N = \beta S_N X^T t$ X se reemplaza por la matriz de características.	Puede operar en espacios no lineales mediante el kernel.	Tiene la misma formulación: $p(t x) = \mathcal{N}(0, K + \sigma^2 I)$ pero \hat{y} es tratado como un valor determinista.	Tiene mayor potencia en modelos no lineal pero su complejidad cubica limita el uso en grandes volúmenes de datos.

Tipo de Regresor	Mínimos Cuadrados	Max Verosimilitud	Max a posteriori	Bayes lineal	Regresión rigida K.	Procesos Gaussianos	Escalabilidad
SGD Regressor	Minimiza la misma función de costo, pero iterativamente usando gradientes por muestra: $w \leftarrow w - \eta \nabla J(w)$	Implementa approx estocástica actualizando parámetros de forma progresiva.	Incluye: Prior Gaussiano // Laplaciano // mixto Adapta según el regresor	Obtiene solo el valor más probable mediante actualizaciones secuenciales.	Trabaja en el espacio original de las características y no usa Kernels.	Se puede considerar como una versión lineal de GP	Tiene eficiencia media. Permite entrenar modelos lineales sobre millones de muestras y variables.
Bayesian Ridge	Incluye regularización e inferencia probabilística: $J(w) = \frac{\beta}{2} \ t - Xw\ ^2 + \frac{\alpha}{2} \ w\ ^2$	La inferencia sobre $p(w)$ es distribución normal. $p(w x, \beta)$ a $p(x w, \beta)$. $p(w \alpha)$.	Adapta una distribución gaussiana sobre $p(w)$.	Version paramétrica del modelo donde α y β se estiman de los datos directamente.	$J(w) = \frac{\beta}{2} \ t - Xw\ ^2 + \frac{\alpha}{2} \ w\ ^2$. Es un ridge probabilístico que obtiene distribuciones de incertidumbre sobre los pesos.	Se considera un GP con espacio de características finito	Combina precisión probabilística con costo razonable. No es tan escalable como SGD Regressor.
Gaussian Process Regressor	Define una distribución sobre funciones. Entrega una familia infinita de curvas con media y varianza. $f(x) \sim g(p(0, K(x, x)))$ sin prior ni incertidumbre	No hay parámetros fijos, se maximiza la verosimilitud marginal sobre las funciones: $p(t x) = N(0, K + \sigma^2 I)$.	Se integra sobre el espacio de funciones completo.	Define una matriz de covarianza $K = xx^T$. Cuantifica la incertidumbre.	Es el Bayesian ridge con infinitas bases implícitas.	La media es igual al resultado del Kernel Ridge, pero sin incertidumbre.	Es el más completo de la familia bayesiana, pero es solo bueno para conjuntos medianos. Tiene incertidumbre real.
Support Vector Machine Regressor	Minimiza la norma de los pesos + penaliza solo los errores mayores a ϵ . Ignora errores pequeños	Tiene max verosimilitud con ruido Laplaciano truncado. Tiene zona de tolerancia	Usa una función de pérdida lineal por tramos	Busca modelar el ruido con una tolerancia fija	Busca minimizar la pérdida g . $f(x) = \sum (a_i - a_i^*) K(x_i, x) + b$.	Es de naturaleza determinista y no se considera gaussiano, solo trabaja con la media.	La versión lineal se usa para muchos datos, pero Kernel SVM solo para datos sets medianos debido a N .
Random Forest Regressor	Ajusta múltiples funciones locales, cada una valida en una región del espacio. Puede aproximar regiones no lineales.	No reduce el error Max Verosimilitud, reduce la varianza por ensemble. Se asume media + ruido gaussiano iid.	No es un estimador MAP clásico. Es un ensemble no paramétrico de MSE por nodo.	La incertidumbre que se reporta suele ser empírica no una varianza bayesiana bien calibrada.	Se puede ver como una regresión Kernel, pues dos puntos pueden tener similitud si caen en la misma hoja y tener un mismo núcleo.	Computo de ideas de promediar puntos similares, pero no hay prior ni posterior probabilística	Es escalable, paralelizable y resistente al ruido, puede estimar incertidumbre por votación entre árboles.

Tipo de Regresor	Min Cuadrados	Max Verosimilitud	Max a posteriori	Bayes lineal	Ridge Kernel	Procesos Gaussianos	Escalabilidad
Gradient Boosting Regressor	Minimiza $\sum L(y, f(x_i))$ por gradiente funcional y tiene solución iterativa. Tiene regularización implícita con el número de árboles.	Minimiza la pérdida cuadrática, que es equivalente a maximizar la verosimilitud. Realiza descenso del gradiente sobre el log-likelihood.	No asume ningún prior ni probabilidad, pero sí regulariza implícitamente. El learning rate actúa como un prior suave.	Ambos son promedios ponderados, pero el GBR es iterativo sobre un modelo flexible.	Tiene linealidad a través de los árboles y suma los árboles que contienen errores, como los kernels.	Ambos construyen modelos de funciones no lineales, pero el GBR lo hace determinísticamente. Si el kernel tiene regiones locales, su media se comporta similar al ajuste por Boosting.	El costo depende del número de árboles M y de muestras N . No es bueno para $N > 1M$, pero es robusto y tiene mejor sensibilidad.
XG Boost	Si se usa $L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$ y $\lambda = 0$, es un GBR, pero con regularización actúa como un Ridge no local.	Minimiza pérdidas equivalentes a log-likelihood negativos, pero la regularización explícita lo convierte en un MAP. Con un modelo no lineal.	Es una implementación determinística del MAP, donde el prior está codificado en los hiperparámetros $(\lambda, \gamma, \alpha)$.	Cada hoja tiene una creencia previa de que los valores deben ser pequeños.	Es una versión discretizada, jerárquica y regularizada del kernel, donde los árboles completan los kernels.	No usa kernels explícitos, pero cada árbol tiene regiones locales donde se agrupan puntos similares. No tiene varianzas ni posterior gaussiano.	Es eficiente y robusto, logra suavidad local y escalabilidad masiva.